

SEP 09 – 11, 2025

CONTAINER
days
CONFERENCE

From IDP to AIDP: Evolving Your Platform for the Machine Learning Age



CONTAINERDAYS CONFERENCE 2025

Max Körbächer



The Convergence of Platform Engineering and AI ?

What is an AI-IDP?

AI-Powered Platforms

Platforms that use AI **to enhance their own capabilities**, providing intelligent assistance, automation, and optimization for development workflows.

AI-Enabling Platforms

Platforms designed to **support the development**, deployment, and management **of AI/ML** workloads, making AI capabilities accessible to developers.

Core Benefits

- Accelerated development cycles
- Improved code quality
- Enhanced developer productivity
- Reduced operational overhead

Bridging the Gap: Developers and AI Engineering Teams

Traditional Developers

Code-first approach

- Small artifacts, predictable CI/CD, auto. testing
- Standard hardware resources
- Scale horizontally

AI/ML Engineers

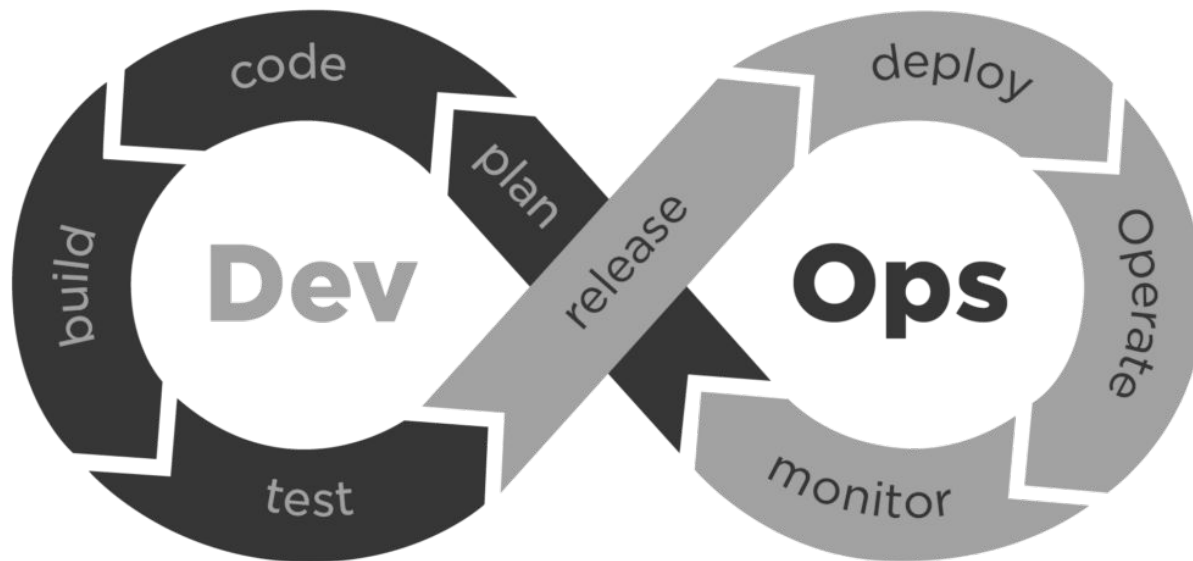
Experiment-first approach

- Iterative experimentation, complex flows
- Large models, massive datasets
- Specialized hardware requirements
- Scale vertically, mix

SEP 09 – 11, 2025

CONTAINER days

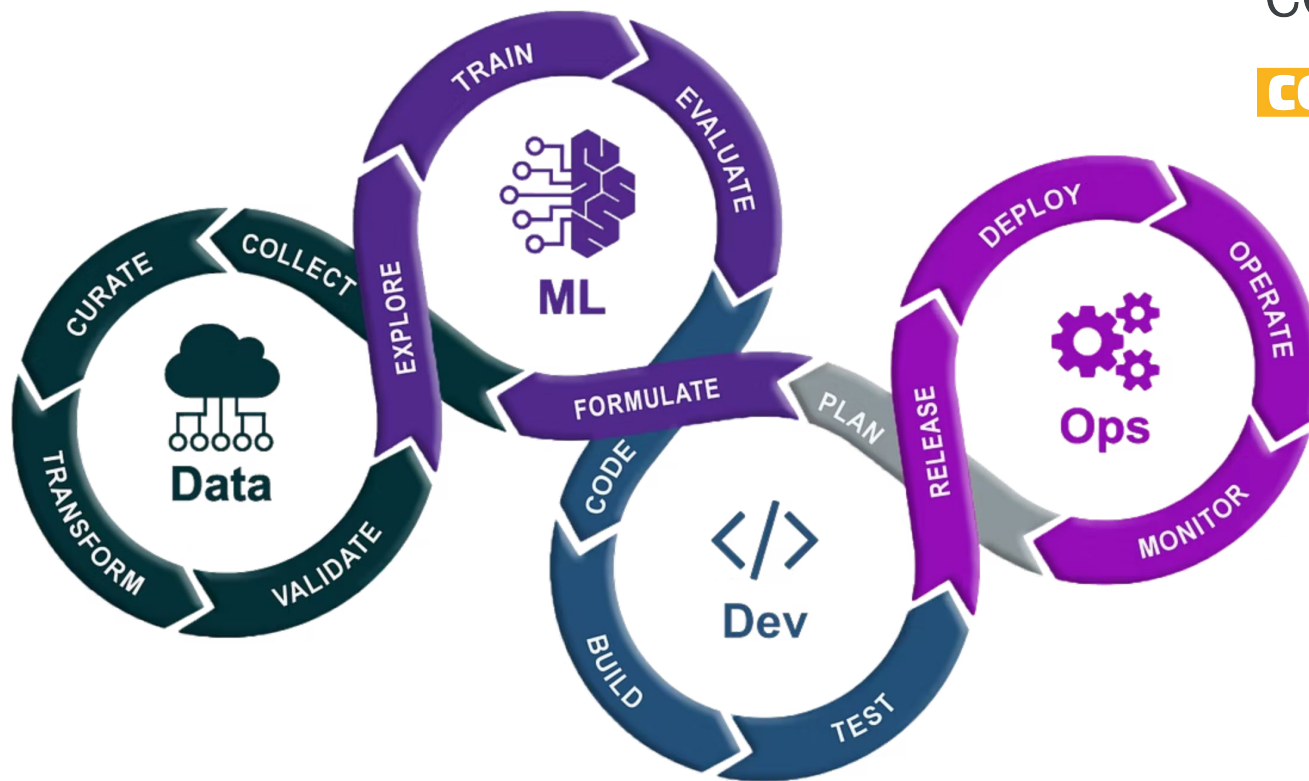
CONFERENCE



SEP 09 – 11, 2025

CONTAINER days

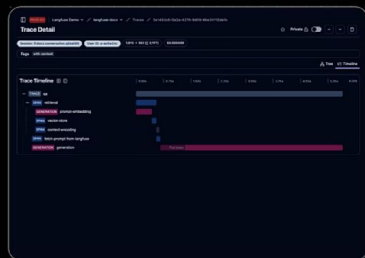
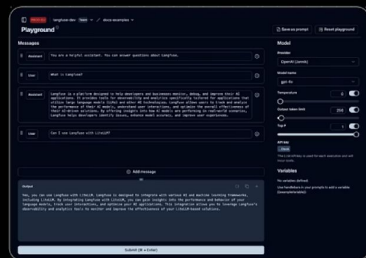
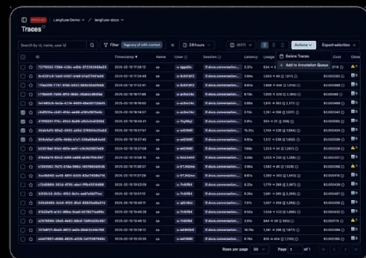
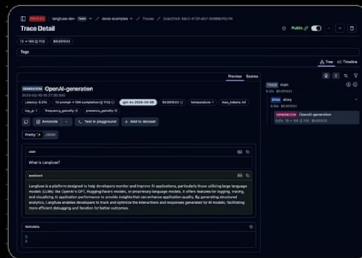
CONFERENCE



Understanding your users

langfuse as example

traces, evals, metrics,
 prompt management,
 and playground...



...to debug and
 improve your LLM
 application together.

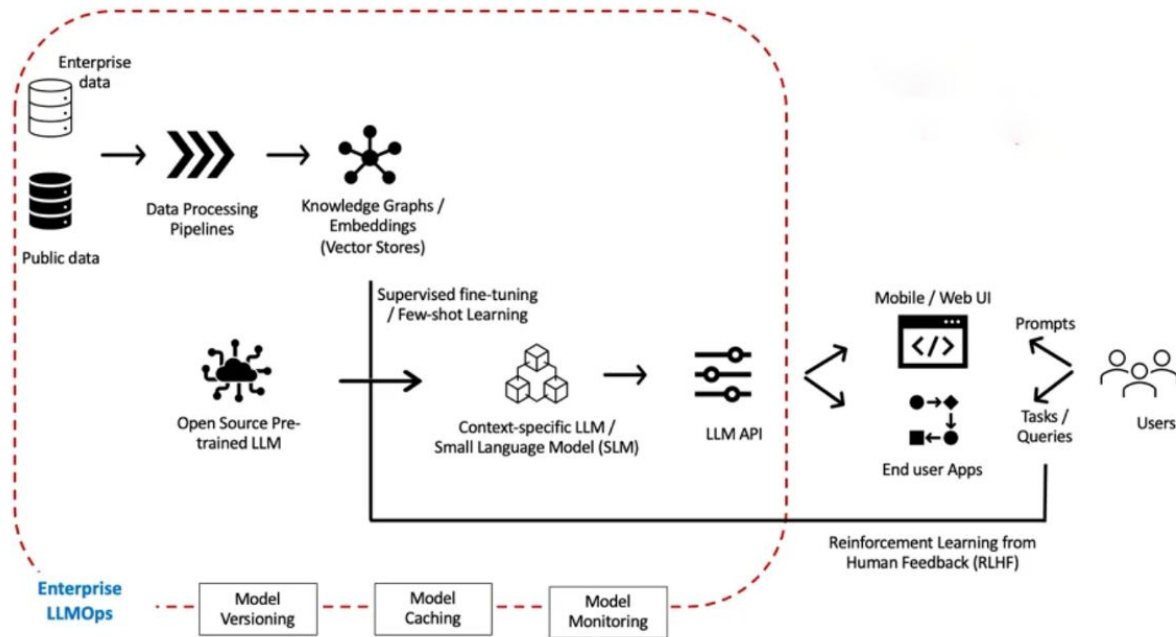
Works with any LLM app and model

SDKs for Python & JS/TS, native integrations for popular libraries and support for OpenTelemetry

Integration overview → Request new integration →

Python SDK	JS/TS SDK	API	OpenTelemetry	OpenAI SDK
LangChain, LangGraph	Llama-Index	CrewAI	LiteLLM	nlp 1.0 AI SDK
Haystack	Instructor	Semantic Kernel	DSPy	Smolagents
Pycardic AI	AutoGen	Amazon Bedrock	Google Vertex/Gemini	Ollama
Flowise	Langflow	Dify	OpenWeb UI	More

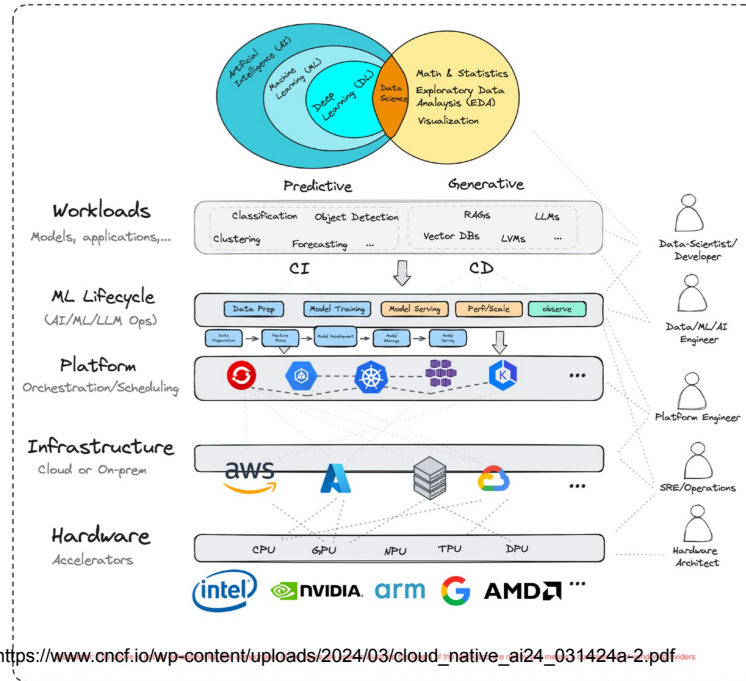
But there is more, right?



Source: Debmalya Biswas

But there is more, right, RIGHT?

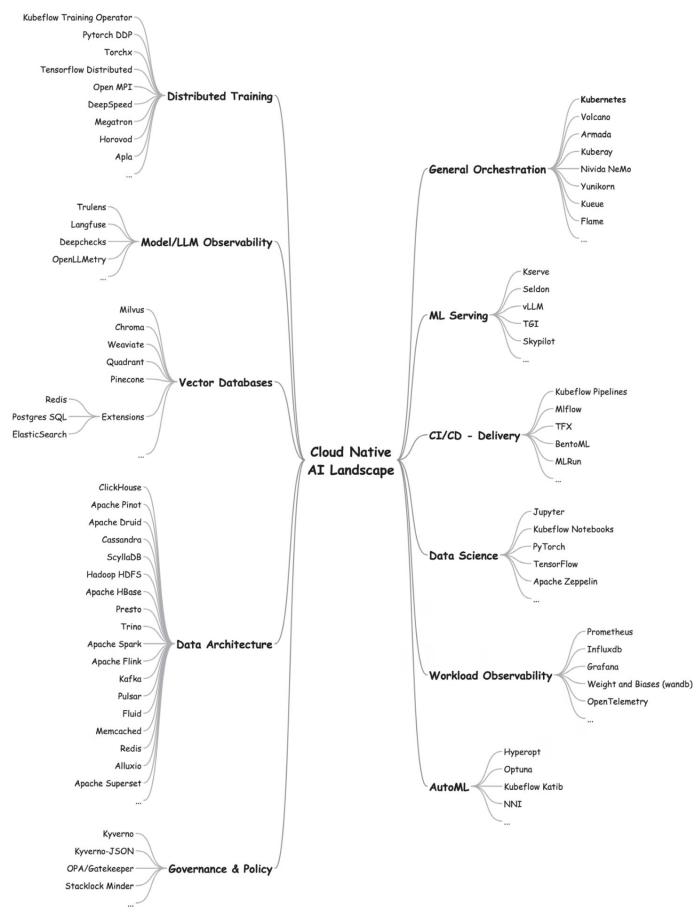
cloud Native AI



Cloud Native Artificial Intelligence Whitepaper: https://www.cncf.io/wp-content/uploads/2024/03/cloud_native_ai24_031424a-2.pdf

Figure 1
Cloud Native AI

But there is more, right, RIGHT; R3!*gT?



SEP 09 – 11, 2025
CONTAINER
days
CONFERENCE

Cloud Native Artificial Intelligence Whitepaper: <https://www.cn>

Figure 4
ML Tool to Task Mind Map

SEP 09 – 11, 2025

CONTAINER
days
CONFERENCE



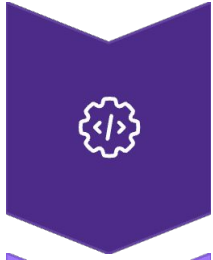
The 3 Elements you have to master first



SEP 09 – 11, 2025

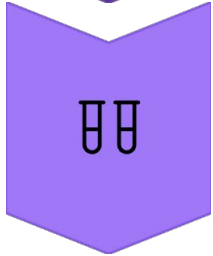
CONTAINER
days

CONFERENCE



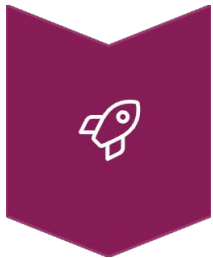
“MLOps” - Pipelines

Orchestrate, manage and scale



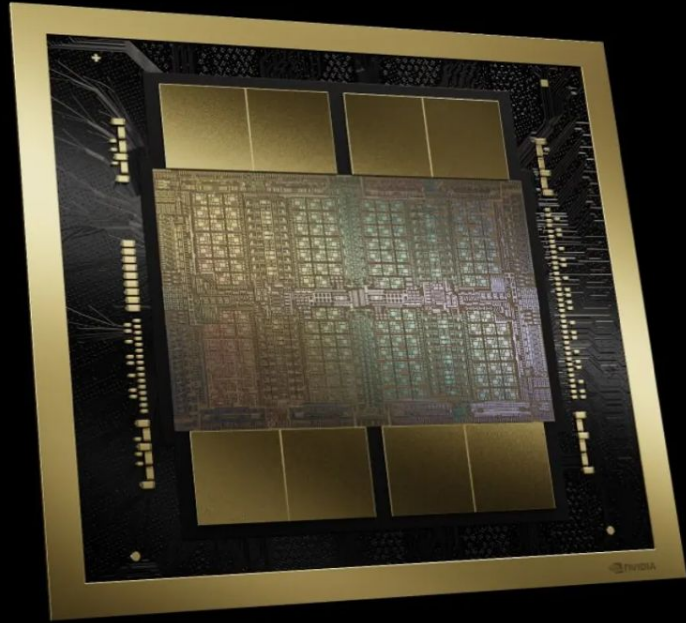
Lifecycle and orchestration of AI workload

That's what we are doing just without AI so far



AI Accelerating Hardware Scheduling & Management

Provide use case oriented hardware



GPU Resource Management

vGPU (Virtual GPU)

Sharing of one GPU with multiple VMs or Containers

MIG (Multi-instance GPU)

Partitioning of a physical GPU for multiple independent workloads

GPU Time-Slicing

Slices GPU time across multiple tasks

CUDA

Software to develop, manage and distribute workload

Physical GPU

```
apiVersion: v1
kind: Pod
metadata:
  name: gpu-example
spec:
  containers:
  - name: gpu-example
    image: nvidia/cuda
    resources:
      limits:
        nvidia.com/gpu: 2
  nodeSelector:
    nvidia.com/gpu.product: A100-PCIE-40GB
    nvidia.com/cuda.runtime: 11.4
    nvidia.com/cuda.driver: 470.161.03
```

MIG

```
apiVersion: v1
kind: Pod
metadata:
  name: gpu-example
spec:
  containers:
  - name: gpu-example
    image: nvidia/cuda
    resources:
      limits:
        nvidia.com/mig-1g.5gb: 1
  nodeSelector:
    nvidia.com/gpu.product: A100-PCIE-40GB
    nvidia.com/cuda.runtime: 11.4
    nvidia.com/cuda.driver: 470.161.03
```

Time Slices with GPU

```
apiVersion: v1
kind: Pod
metadata:
  name: gpu-example
spec:
  containers:
  - name: gpu-example
    image: nvidia/cuda
    resources:
      limits:
        nvidia.com/gpu.shared: 1
  nodeSelector:
    nvidia.com/gpu.product: A100-PCIE-40GB
    nvidia.com/cuda.runtime: 11.4
    nvidia.com/cuda.driver: 470.161.03
```

```
version: v1
sharing: k8s-device-plugin
timeSlicing: config file
resources:
  - name: nvidia.com/gpu
    replicas: 10
  ...
```

Time Slices with MIG

```
apiVersion: v1
kind: Pod
metadata:
  name: gpu-example
spec:
  containers:
  - name: gpu-example
    image: nvidia/cuda
    resources:
      limits:
        nvidia.com/mig-1g.5gb.shared: 1
  nodeSelector:
    nvidia.com/gpu.product: A100-PCIE-40GB
    nvidia.com/cuda.runtime: 11.4
    nvidia.com/cuda.driver: 470.161.03
```

```
version: v1
sharing: k8s-device-plugin
timeSlicing: config file
resources:
  - name: nvidia.com/gpu
    replicas: 10
  - name: nvidia.com/mig-1g.5gb
    replicas: 10
  ...
```

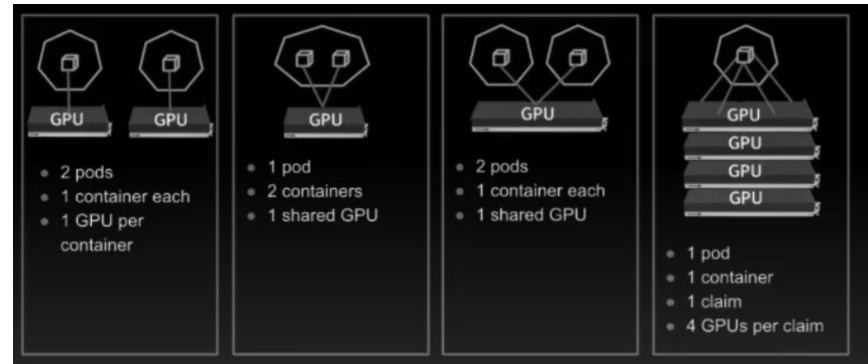
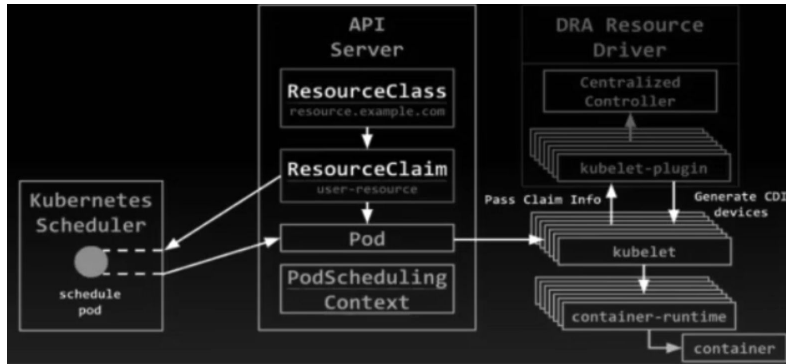
Limitations with this approach

- Heterogeneous GPU Support**
No native support for integrating more than one GPU type per node, restricting hardware configurations.
- Complex Constraint Handling**
Lack of robust mechanisms for providing complex constraints when requesting a GPU, making advanced scheduling difficult.
- Oversubscription Control**
Limited control over how oversubscribed GPUs are shared between jobs, potentially leading to unfair resource distribution.
- MPS Integration Challenges**
Awkward and overly-burdensome support for Multi-Process Service (MPS), adding operational overhead.
- Dynamic MIG Provisioning**
Inability to dynamically provision Multi-Instance GPU (MIG) devices based on incoming requests, reducing resource elasticity.
- Driver Selection Flexibility**
No built-in capability to dynamically choose between NVIDIA and other drivers (e.g., fio) on a per-GPU basis.

DRA - Dynamic Resource Allocation

Allocate **cross-node resources** in Kubernetes

Explicitly **share, partition, and reconfigure** devices **on-the-fly** based on user requests



DRA - Dynamic Resource Allocation

```

---
apiVersion: resource.k8s.io/v1alpha2
kind: ResourceClaimTemplate
metadata:
  name: unique-gpu
spec:
  spec:
    resourceName: gpu.nvidia.com
---
apiVersion: v1
kind: Pod
metadata:
  name: gpu-example
spec:
  containers:
  - name: ctr
    image: nvidia/cuda
    command: ["nvidia-smi" "-L"]
    resources:
      claims:
      - name: gpu0
      - name: gpu1
  resourceClaims:
  - name: gpu0
    source:
      resourceClaimTemplateName: unique-gpu
  - name: gpu1
    source:
      resourceClaimTemplateName: unique-gpu

```

Associated with the DRA Driver and installed by the cluster admin

```

apiVersion: v1
kind: Pod
metadata:
  name: gpu-example
spec:
  containers:
  - name: ctr0
    resources:
      claims:
      - name: gpu
  - name: ctr1
    resources:
      claims:
      - name: gpu
  resourceClaims:
  - name: gpu
    source:
      resourceClaimName: unique-gpu

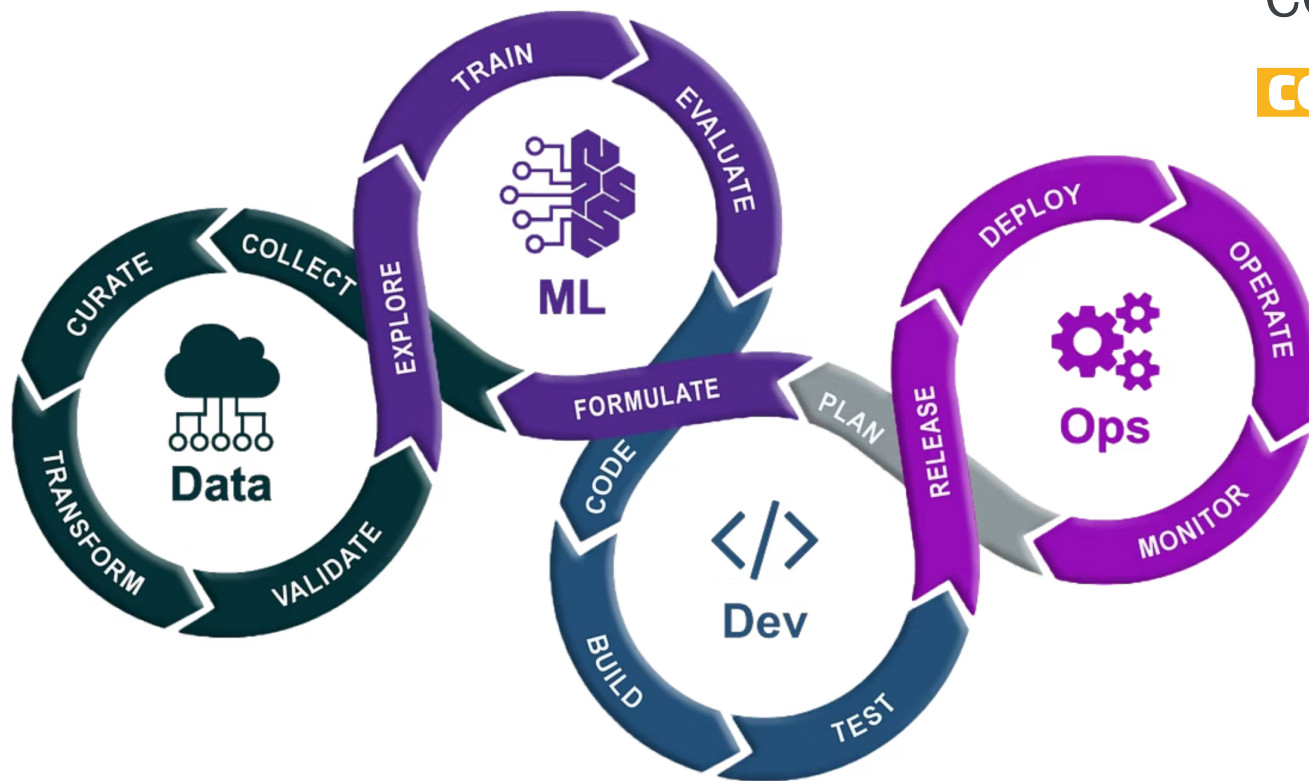
```

Shared access to same underlying GPU

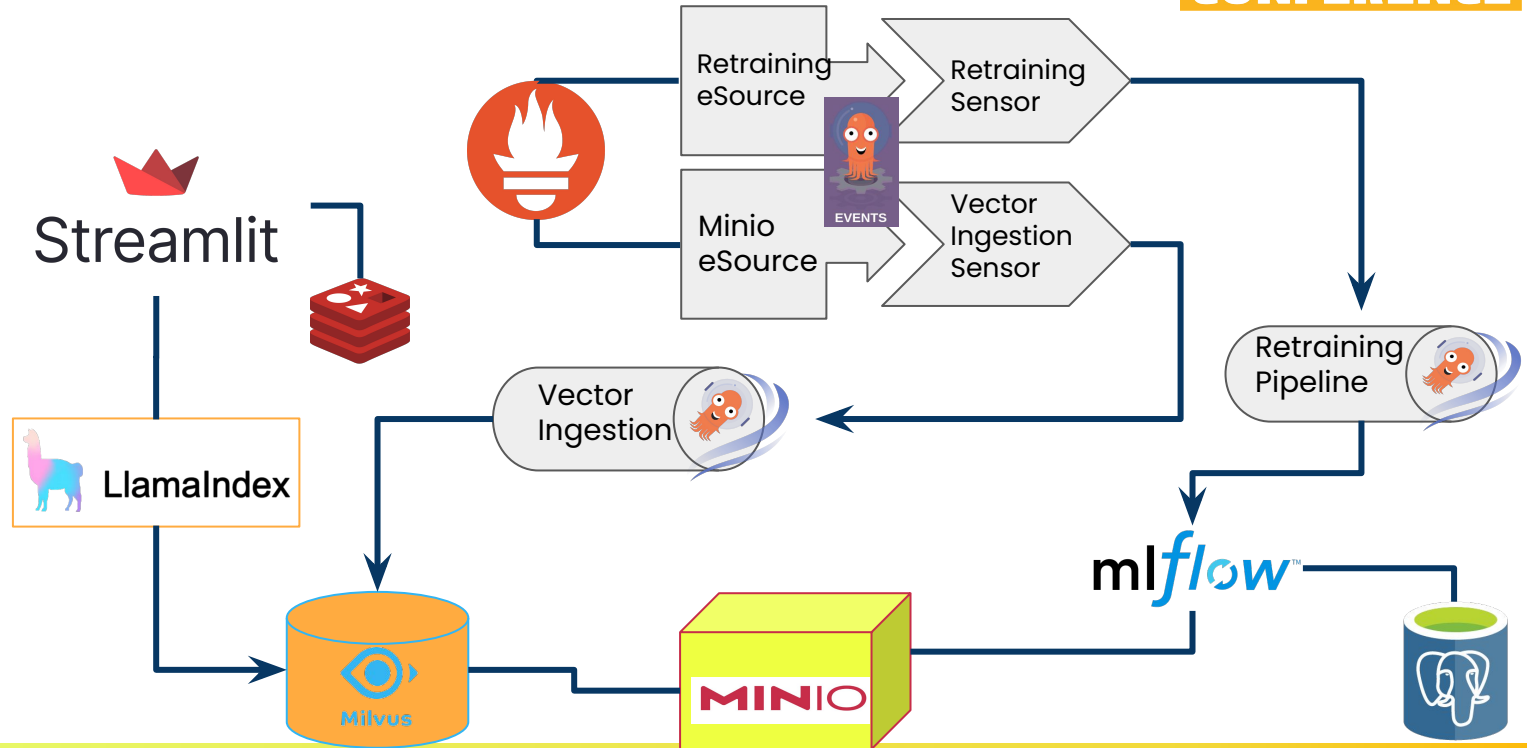
SEP 09 – 11, 2025

CONTAINER days

CONFERENCE



ML/LLMOps example Architecture



Kubeflow Ecosystem

AI Ecosystem

JupyterLab

VSCode

RStudio

PyTorch

HuggingFace

TensorFlow

DeepSpeed

XGBoost

Megatron-LM

Horovod

Scikit-Learn

MPI

Optuna

Hyperopt

others...

Kubeflow Projects

Kubeflow Projects

Kubeflow
Spark Operator

Kubeflow
Trainer

Kubeflow
Katib

KServe

Kubeflow
Notebooks

Kubeflow
Pipelines

Kubeflow
Dashboard

Kubeflow
Model Registry

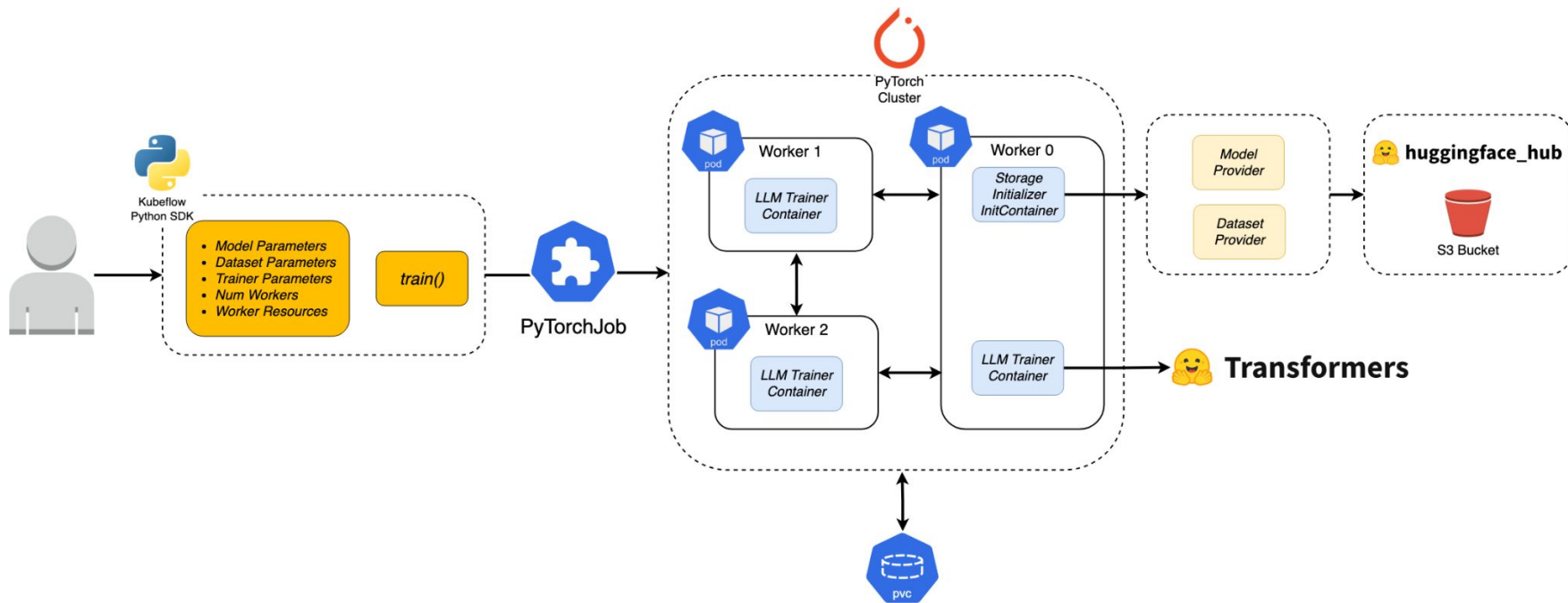
External Add-Ons

Feast

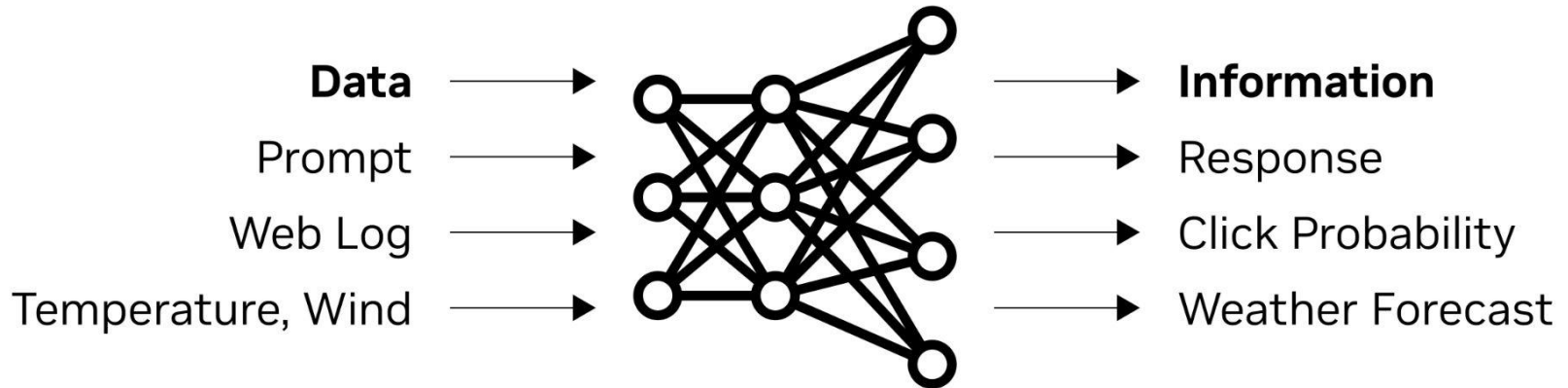
Elyra

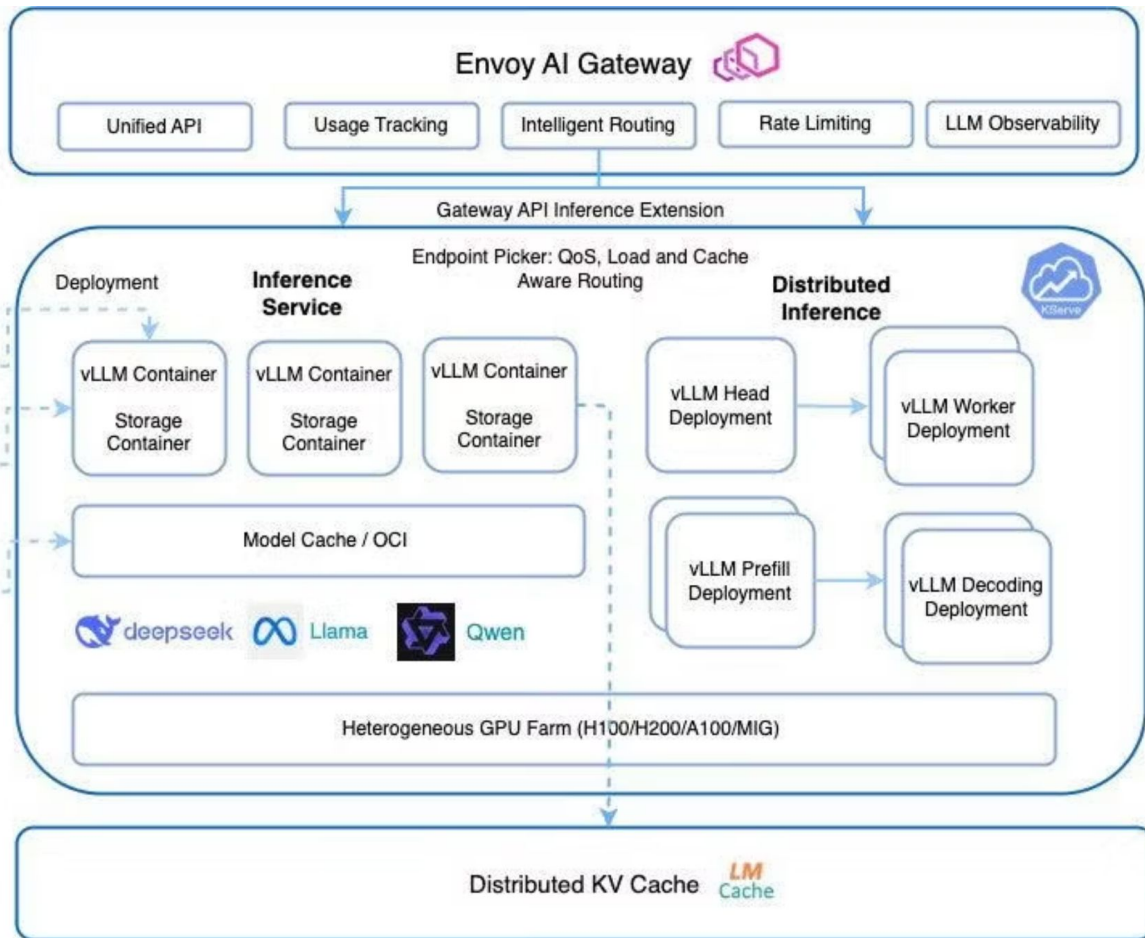
others...

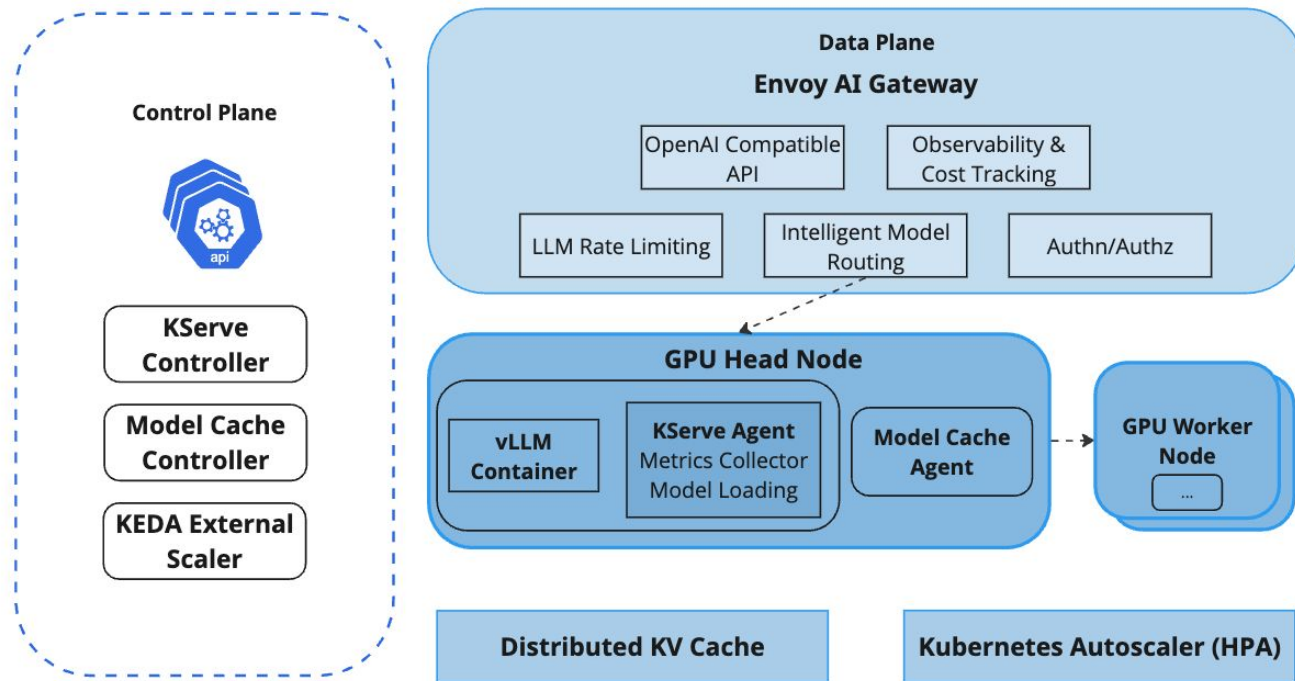
Fine-tuning with kflow



General Inference







SEP 09 – 11, 2025

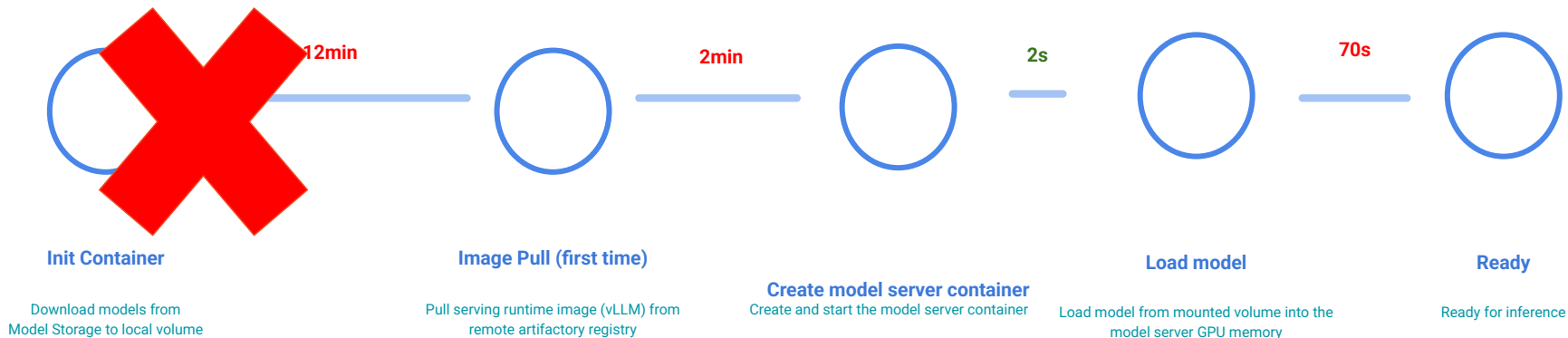
CONTAINER
days
CONFERENCE



3 fast steps to optimize your platform



Model Caching



Prompt Caching

KVCache

Grows too fast! 

LMCache

Reuses common prefixes.

 Faster inference!

KV Cache Size Calculator

Select LLM Model:

deepseek-ai/DeepSeek-V3

Select data type:

float16

Enter Number of Tokens:

10000

Calculate KV Cache Size

KV Cache Size: 16.2888 GB

Calculation Details:

Selected Model: deepseek-ai/DeepSeek-V3
Hidden Size: 7168
Number of Attention Heads: 128
Number of Hidden Layers: 61
Number of Key-Value Heads: 128
Head Size: 56 (Hidden Size / Attention Heads)
Data Type Size: 2 bytes
Total Elements: $2 \times 61 \times 10000 \times 128 \times 56 = 8744960000$
Total Bytes: $8744960000 \times 2 = 17489920000$ bytes
KV Cache Size: $17489920000 / (1024^2) \approx 16.2888$ GB

KV Cache Size Calculator

Select LLM Model:

deepseek-ai/DeepSeek-V3

Select data type:

float16

Enter Number of Tokens:

1000000

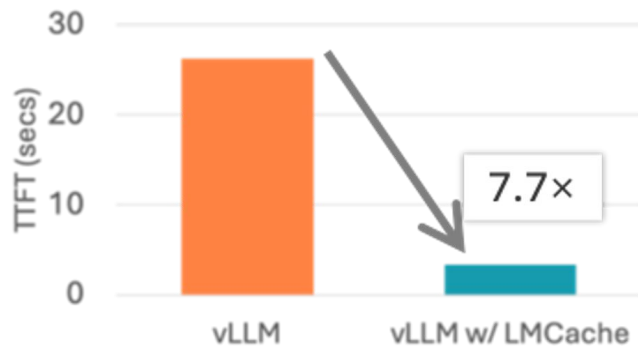
Calculate KV Cache Size

KV Cache Size: 1628.8757 GB

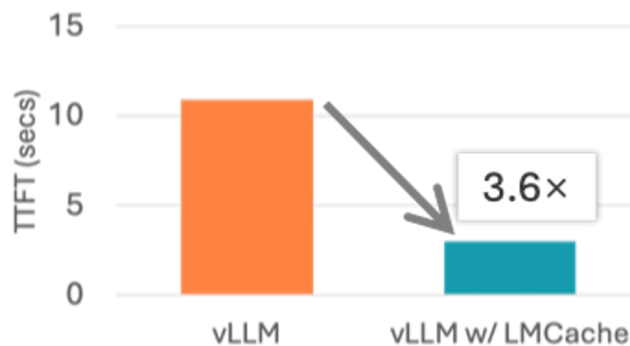
Calculation Details:

Selected Model: deepseek-ai/DeepSeek-V3
Hidden Size: 7168
Number of Attention Heads: 128
Number of Hidden Layers: 61
Number of Key-Value Heads: 128
Head Size: 56 (Hidden Size / Attention Heads)
Data Type Size: 2 bytes
Total Elements: $2 \times 61 \times 1000000 \times 128 \times 56 = 874496000000$
Total Bytes: $874496000000 \times 2 = 1748992000000$ bytes
KV Cache Size: $1748992000000 / (1024^2) \approx 1628.8757$ GB

Prompt Caching



Use case 1: **long context**
Context length: **25K** tokens
(Llama 70B on A40)



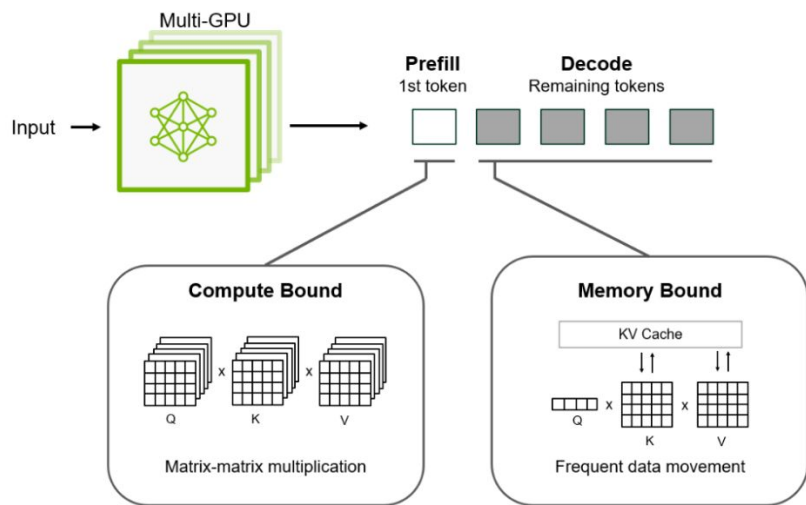
Use case 2: **RAG**
Retrieved chunks: **4 x 2K** tokens
(Llama 70B on A40)

LMCache drastically reduces the prefill delay (TTFT) by reusing KV caches

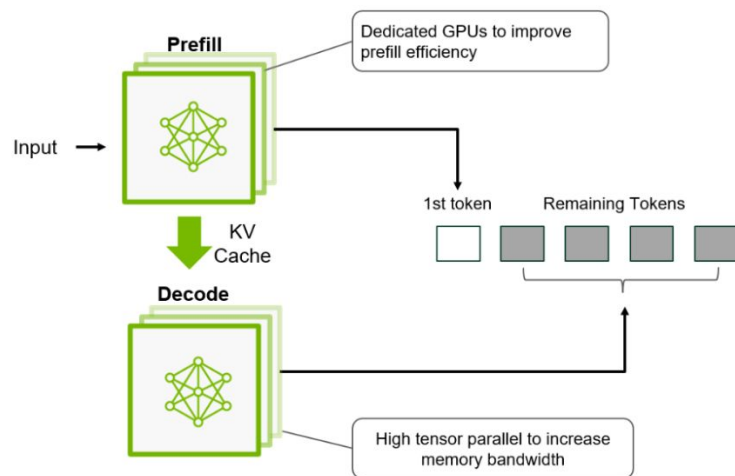
- > 3-10x delay savings
- > GPU cycle reduction in use cases like multi-round QA and RAG

Disaggregated Serving

Traditional Serving



Disaggregated Serving



<https://developer.nvidia.com/blog/introducing-nvidia-dynamo-a-low-latency-distributed-inference-framework-for-scaling-reasoning-ai-models/>

Disaggregated Serving

Tokens Per Second Per GPU
(Higher is Better)



SEP 09 – 11, 2025

CONTAINER
days
CONFERENCE

Closing 



SEP 09 – 11, 2025

CONTAINER
days

CONFERENCE

Measuring the Impact of Your AI-Ready IDP

AI Serving Signals



Token Consumption

Daily processing volume



Guardrail Success

Protection execution rate



PII Leakage

Sensitive data exposure rate

AI Performance Metrics



Inference Latency

Average response time



Availability

System uptime for AI services



Throughput Gain

vs. non-optimized
deployment

Platforms for LLMs Enable ISO 42001 & EU AI Act

AI-ready Internal Developer Platforms unlock compliance with emerging AI governance standards.

2023

Standard Release

First international standard for AI management systems.

42001

Certification

Demonstrates commitment to responsible AI practices.

100%

Governance

Complete oversight of AI assets and processes.

SEP 09 – 11, 2025
CONTAINER
days
CONFERENCE

Common Pitfalls

Over-Engineering from Day One

Underestimating Resource Requirements

Neglecting Model Governance

Poor Developer Adoption

Key Takeaways

Balance Control & Flexibility

Create standardized paths that make the right way the easy way, while allowing for customization when necessary.

Start Simple, Iterate Often

Begin with core capabilities that unblock key workflows. Add sophistication based on actual user needs and feedback.

Developer Experience Drives Adoption

The most technically impressive platform is worthless if developers avoid using it. Invest in intuitive interfaces and documentation.

Efficient Resource Management is Critical

GPU utilization directly impacts both cost-effectiveness and user satisfaction. Optimization pays dividends at scale.